

TREE ENSEMBLES WITH RULE STRUCTURED HORSESHOE REGULARIZATION

MALTE NALENZ AND MATTIAS VILLANI

ABSTRACT. We propose a new Bayesian model for flexible nonlinear regression and classification using tree ensembles. The model is based on the RuleFit approach in [9] where rules from decision trees and linear terms are used in a L1-regularized regression. We modify RuleFit by replacing the L1-regularization by a horseshoe prior, which is well known to give aggressive shrinkage of noise predictor while leaving the important signal essentially untouched. This is especially important when a large number of rules are used as predictors as many of them only contribute noise. Our horseshoe prior has an additional hierarchical layer that applies more shrinkage a priori to rules with a large number of splits, and to rules that are only satisfied by a few observations. The aggressive noise shrinkage of our prior also makes it possible to complement the rules from boosting in [9] with an additional set of trees from random forest, which brings a desirable diversity to the ensemble. We sample from the posterior distribution using a very efficient and easily implemented Gibbs sampler. The new model is shown to outperform state-of-the-art methods like RuleFit, BART and random forest on 16 datasets. The model and its interpretation is demonstrated on the well known Boston housing data, and on gene expression data for cancer classification. The posterior sampling, prediction and graphical tools for interpreting the model results are implemented in a publicly available R package.

1. INTRODUCTION

Learning and prediction when the mapping between input and outputs is potentially nonlinear and observed in noise remains a major challenge. Given a set of N training observations $(\mathbf{x}, y)_i, i = 1, \dots, N$, we are interested in learning or approximating an unknown function f observed in additive Gaussian noise

$$y = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2),$$

and to use the model for prediction. A popular approach is to use a learning ensemble

$$f(\mathbf{x}) = \sum_{l=1}^m \alpha_l f_l(\mathbf{x}),$$

where $f_l(\mathbf{x})$ is a basis function (also called a weak learner in the machine learning literature) for a subset of the predictors. A variety of basis functions f_l have been proposed in the last decades, and we will here focus on decision rules. Decision rules are defined by simple if-else statements and therefore highly interpretable by humans. Finding a set of optimal rules is NP hard [9], and most practical algorithms therefore use a greedy learning procedure. Among the most powerful are divide and conquer approaches [5, 10] and boosting [17, 6].

Nalenz: *Division of Statistics, Dept. of Computer and Information Science, SE-581 83 Linköping, Sweden. E-mail: malte.nalenz@helmholtz-muenchen.de.* Villani: *Division of Statistics and Machine Learning, Dept. of Computer and Information Science, SE-581 83 Linköping, Sweden. E-mail: mattias.villani@liu.se.*

A fundamental new way to learn decision rules is introduced in Popescu and Friedman [9] in their RuleFit approach. RuleFit is estimated by a two-step procedure. The *rule generation* step extracts decision rules from an ensemble of trees trained with gradient boosting. The second *regularization* step learns the weights α_l for the generated rules via L1-regularized (Lasso) regression, along with weights on linear terms included in the model. RuleFit has been successfully applied in particle physics, in medical informatics and in life sciences. Our paper makes the following contributions to improve and enhance RuleFit.

First, we replace the L1-regularization [20] in RuleFit by a generalized horseshoe regularization prior [3] tailored specifically to covariates from a rule generation step. L1-regularization is computationally attractive, but has the well known drawback of also shrinking the effect of the important covariates. This is especially problematic here since the number of rules from the rule generation step can be very large while potentially only a small subset is necessary to explain the variation in the response. Another consequence of the overshrinkage effect of the L1-regularization is that it is hard to choose an optimal number of rules; increasing the number of rules affects the shrinkage properties of the Lasso. This makes it very hard to determine the number of rules a priori, and one has to resort to cross-validation, thereby mitigating the computational advantage of the Lasso. A horseshoe prior is especially attractive for rule learning since it shrinks uninformative predictors aggressively while leaving important ones essentially untouched. Inspired by the prior distribution on the tree depth in BART [4], we design a generalized horseshoe prior that shrinks overly complicated and specific rules more heavily, thereby mitigating problems with overfitting. This is diametrically opposed to RuleFit, and to BART and boosting, who all combine a myriad of rules into a collective where single rules only play a very small part.

Second, we complement the tree ensemble from gradient boosting [7] in RuleFit [9] with an additional set of trees generated with Random Forest. The error-correcting nature of boosting makes the rules highly dependent on each other. Trees from Random Forest [1] are much more random and adding them to rules from boosting therefore brings a beneficial diversity to the tree ensemble. This agrees well with the aggressive horseshoe regularization since it automatically prunes away noise trees and avoids overfitting.

Third, an advantage of our approach compared to many other flexible regression and classification models is that predictions from our model are based on a relatively small set of interpretable decision rules. The possibility to include linear terms also simplifies interpretation since it avoids a common problem with decision trees that linear relationships need to be approximated with a large number of rules. To further aid in the interpretation of the model and its predictions, we also propose graphical tools for analyzing the model output. We also experiment with post-processing methods for additional pruning of rules to simplify the interpretation even further using the method in [12].

We call the resulting two-step procedure with mixed rule generation followed by generalized rule structured horseshoe regularization the *HorseRule* model. We show that HorseRule's ability to keep the important rules and aggressively removing unimportant noise rules leads to both great predictive performance and high interpretability.

The structure of the paper is as follows. Section 2 describes the decision rule generation method in Hs-RuleFit. Section 3 presents the horseshoe regularization prior and the MCMC algorithm for posterior inference. Section 4 illustrates aspects of the approach on simulated data and evaluates and compares the predictive performance of Hs-RuleFit to several main competing methods on a wide variety of data sets. Section 5 concludes.

2. DECISION RULE GENERATION

This section describes the *rule generation step* of HorseRule, which complements the rules from gradient boosting in Friedman and Popescu [9] with rules from Random Forest with completely different properties.

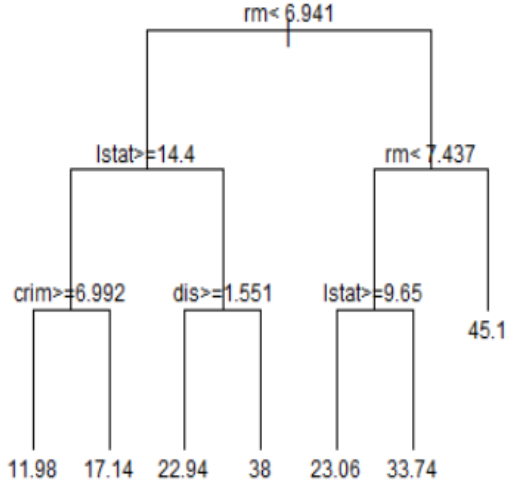


FIGURE 2.1. Decision Tree for the Boston Housing data.

Rules	Conditions
r_1	$RM \geq 6.94$
r_2	$RM < 6.94$
r_3	$RM < 6.94 \ \& \ LSTAT < 14.4$
r_4	$RM < 6.94 \ \& \ LSTAT \geq 14.4$
r_5	$RM < 6.94 \ \& \ LSTAT < 14.4 \ \& \ CRIM < 6.9$
r_6	$RM < 6.94 \ \& \ LSTAT < 14.4 \ \& \ CRIM \geq 6.9$
r_7	$RM \geq 6.94 \ \& \ LSTAT < 14.4 \ \& \ DIS < 1.5$
r_8	$RM \geq 6.94 \ \& \ LSTAT < 14.4 \ \& \ DIS \geq 1.5$
r_9	$6.94 \leq RM < 7.45$
r_{10}	$6.94 \leq RM < 7.45$
r_{11}	$6.94 \leq RM < 7.45 \ \& \ LSTAT < 9.7$
r_{12}	$6.94 \leq rm < 7.45 \ \& \ LSTAT \geq 9.7$

TABLE 2.1. Corresponding rules, defining the Decision Tree.

2.1. Decision Rules. Let S_k denote the set of possible values of the covariate x_k and let $s_{k,m} \subseteq S_k$ denote a specific subset. A decision rule can then be written as

$$(2.1) \quad r_m(x) = \prod_{k \in Q_m} I(x_k \in s_{k,m}),$$

where $I(x)$ is the indicator function and Q_m is the set of variables used in defining the m th rule. A decision rule $r_m \in \{0, 1\}$ takes the value 1 if all of its $|Q_m|$ conditions are fulfilled and 0 otherwise. For orderable covariates $s_{k,m}$ will be an interval or a disjoint union of intervals, while for categorical covariates $s_{k,m}$ are explicitly enumerated. There is a long tradition in machine learning to use decision rules as weak learners. Most algorithms learn decision rules directly from data, such as in [5, 6]. RuleFit exploits the fact that decision trees can be seen as a set of decision rules. In a first step a tree ensemble is generated, which is then decomposed into its defining decision rules. Several efficient (greedy) algorithmic implementations are available for constructing the tree ensembles. The generated rules typically correspond to interesting subspaces with great predictive power. Each node in a decision tree is defined by a decision rule. Figure 1.1 shows an example tree for the Boston Housing data set and Table

1 its corresponding decision rules. Using Equation (2.1) for example r_{11} can be expressed as

$$r_{11}(x) = \prod_{k \in Q_{11}} I(x_k \in s_{k,11}) = I(6.94 \leq RM < 7.45)I(LSTAT < 9.7).$$

The m th tree consists of $2(u_m - 1)$ rules, where u_m denotes the number of terminal nodes. Therefore $\sum_{m=1}^M 2(u_m - 1)$ rules can be extracted from a tree ensemble of size M .

2.2. Collinearity structure of trees. The generated rules will be combined in a linear model and collinearity is a concern. For example, the two first child nodes in each tree are perfectly negative correlated. Furthermore, each parent node is perfectly collinear with its two child nodes, as it is their union. One common way to deal with the collinearity problem is to include the terminal nodes only. This approach also reduces the number of rules and therefore simplifies computations. We have nevertheless choose to consider all possible rules including also non-terminal ones, but to randomly select one of the two child nodes at each split. The reason for also including non-terminal nodes is three-fold. First, even though each parent node in a tree can be reconstructed as a linear combination of terminal nodes, when using regularization this equivalence no longer holds. Second, our complexity penalizing prior in Section 3.3 is partly based on the number of splits to measure the complexity of a rule, and will therefore shrink the several complex child nodes needed to approximate a simpler parent node. Third, the interpretation of the model is substantially simplified if the model can select a simple parent node instead of many complex child nodes.

2.3. Generating an informative and diverse rule ensemble. Any tree method can be used to generate decision rules. Motivated by the experiments in [8], Rulefit uses gradient boosting for rule generation [9]. Gradient boosting [7] fits each tree iteratively on the pseudo residuals of the current ensemble in an attempt to correct mistakes made by the previous ensemble. This procedure introduces a lot of dependence between the members of the ensemble, and many of the produced rules tend to be informative only when combined to an ensemble. It might therefore not be possible to remove a lot of the decision rules without destroying this dependency structure.

Random Forest on the other hand generates trees independently from all previous trees [1]. Each tree tries to find the individually best partitioning, given a random subset of observations and covariates. Random forest will often generate rules with very similar splits, and the random selection of covariates forces it to often generate decision rules based on uninformative predictors. Random Forest will therefore produce more redundant and uninformative rules compared to gradient boosting, but the generated rules with strong predictive power are not as dependent on rest of the ensemble.

Since the rules from boosting and random forest are very different in nature, it makes sense to use both types of rules to exploit both methods' advantages. This naturally leads to a larger number of candidate rules, but the generalized horseshoe shrinkage proposed in Section 3.2 and 3.3 can very effectively handle redundant rules.

The tuning parameters used in tree generator determine the resulting decision rules. The most impactful is the tree-depth, controlling the complexity of the resulting rules. We follow

[9] with setting the depth of tree m to

$$(2.2) \quad td_m = 2 + \lfloor \varphi \rfloor$$

where $\lfloor x \rfloor$ is the largest integer less or equal than x and φ is a random variable following the exponential distribution with mean $L - 2$. Setting $L = 2$ will produce only tree stumps consisting of one split. With this indirect specification the forest is composed of trees of varying depth, which allows the model to be more adaptive to the data and makes the choice of a suitable tree depth less important.

Another important parameter is the minimum number of observations in a node n_{min} . A too small n_{min} gives very specific rules and the model is likely to capture spurious relationships. Using $n_{min} = N^{\frac{1}{3}}$ as a default setting has worked well in our experiments, but if prior information about reasonable sizes of subgroups in the data is available the parameter can be adjusted accordingly. Another choice is to determine n_{min} by cross validation. In the following all other tuning parameters, e.g the shrinkage parameter in gradient boosting or the number of splitting covariates in the random forest, are set to their recommended standard choices implemented in the R-packages *randomForest* and *gbm*.

3. ENSEMBLES AND RULE BASED HORSESHOE REGULARIZATION

This section discusses the *regularization step* of HorseRule and present a new horseshoe shrinkage prior tailored specifically for covariates in the form of decision rules.

3.1. The ensemble. Once a suitable set of decision rules is generated, they can be combined in a linear regression model of the form

$$y = \alpha_0 + \sum_{l=1}^m \alpha_l r_l(\mathbf{x}) + \epsilon.$$

As $r_i(\mathbf{x}) \in \{0, 1\}$ they already have the form of dummy variables and can be directly included in the regression model. A simple but important extension is to also include linear terms

$$(3.1) \quad y = \alpha_0 + \sum_{j=1}^p \beta_j x_j + \sum_{l=1}^m \alpha_l r_l(\mathbf{x}) + \epsilon.$$

This extension addresses the difficulty of rule and tree based methods to approximate linear effects. Splines, polynomials, time effects, spatial effects or random effects are straightforward extensions of Equation (3.1).

Friedman and Popescu [9] do not standardize the decision rules, which puts a higher penalty on decision rules with a smaller scale. To avoid this behavior, we scale the predictors to have zero mean and unit variance.

3.2. Bayesian regularization through the horseshoe prior. A large set of candidate decision rules is usually necessary to have a high enough chance of finding good decision rules. The model in (3.1) will therefore always be high dimensional and often $p + m > n$. Many of the rules will be uninformative and correlated with each other. Regularization is therefore a necessity.

RuleFit [9] uses L1-regularized estimates, which corresponds to an a posteriori mode estimator under a double exponential prior in a Bayesian framework [20]. As discussed in the Introduction, the global shrinkage effect of L1-regularization can be problematic for rule covariates. L1-regularization is well known to lead to both shrinkage and variable selection. There now exist implementations of RuleFit that use the elastic net instead of L1-Regularization, which can lead to improved predictive performance [21], however elastic net still only uses one global shrinkage parameter.

Another common Bayesian variable selection approach is based on the spike-and-slab prior [11, 19]

$$(3.2) \quad \beta_j \sim w \cdot N(\beta_j; 0, \lambda^2) + (1 - w) \cdot \delta_0,$$

where δ_0 is the Dirac point mass function, $N(\beta_j; 0, \lambda^2)$ is the normal density with zero mean and variance λ^2 , and w is the prior inclusion probability of predictor x_j . Discrete mixture priors enjoy attractive theoretical properties, but need to explore a model space of size $2^{(p+m)}$, which can be problematic when either p or m are large. The horseshoe prior by [2, 3] mimics the behavior of the spike-and-slab but is computationally more attractive. The original horseshoe prior for linear regression is of the form

$$(3.3) \quad y|\mathbf{X}, \boldsymbol{\beta}, \sigma^2 \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n),$$

$$(3.4) \quad \beta_j|\lambda_j, \tau^2, \sigma^2 \sim \mathcal{N}(0, \lambda_j \tau^2 \sigma^2),$$

$$(3.5) \quad \sigma^2 \sim \sigma^{-2} d\sigma^2,$$

$$(3.6) \quad \lambda_j \sim \mathcal{C}^+(0, 1),$$

$$(3.7) \quad \tau \sim \mathcal{C}^+(0, 1).$$

The horseshoe shrinkage for β_j is determined by a local shrinkage parameter $\lambda_j > 0$ and a global shrinkage parameter $\tau > 0$. This is important since it allows aggressive shrinking of noise covariates through small values of τ , while allowing individual signals to have large coefficients through large λ_j . Carvalho et al. [3] show that the horseshoe is better at recovering signals than the Lasso, and the models obtained from the horseshoe are shown to be almost indistinguishable from the ones obtained by a well defined spike-and-slab prior.

3.3. Horseshoe regularization with rule structure. The original horseshoe assigns the same prior distribution to all regression coefficients, regardless of the rule's complexity (number of splits in the tree) and the specificity (number of data points that fulfill the rule). Similar to the tree structure prior in BART, we therefore modify the horseshoe prior to express the prior belief that rules with high length (many conditions) are less likely to reflect a true mechanism. In addition, we also add the prior information that very specific rules that are satisfied by only a few data points are also improbable a priori. These two sources of prior information are incorporated by extending the prior on λ_j to

$$\lambda_j \sim \mathcal{C}^+(0, A_j),$$

with

$$(3.8) \quad A_j = \frac{(2 \cdot \min(1 - s(r_j), s(r_j)))^\mu}{(l(r_j))^\eta},$$

where $l(r_j)$ denotes the length of rule j defined as its number of conditions, and $s(r_l) \in (0, 1)$ denotes the rule support $s(r_j) = N^{-1} \sum_{i=1}^N r_j(\mathbf{x}_i)$. The hyperparameter μ controls the strength of our belief to prefer general rules that cover a lot of observations and η determines how strongly we prefer simple rules. The response y should be scaled when using the rule structure prior since the scale of β depends on the scale of y .

The rule structure for A_j in equation 3.8 is designed such that $A_j = 1$ for rules with support 0.5 and length 1, as the ideal. Since $\lim_{\mu \rightarrow 0, \eta \rightarrow 0} A_j = 1$, our rule structure prior approaches the standard horseshoe prior for small μ and η . The rule structure prior gives a gentle push towards simple and general rules, but its Cauchy tails put considerable probability mass on non-zero values even for very small A_j ; the data can therefore overwhelm the prior and keep a complex and specific rule if needed.

A model with many complex specific rules may drive out linear terms from the model, thereby creating an unnecessarily complicated model. Setting the parameters μ and η to values larger than 0 will give linear effects a higher chance of being chosen a priori. The hyperparameters μ and η can be chosen guided by theoretical knowledge about what kind of rules and linear effects are reasonable for a problem by hand, or determined via cross validation. As a default choice $(\mu, \eta) = (1, 2)$ worked well in our experiments, penalizing rule complexity heavily and low rule support moderately.

3.4. Posterior inference via Gibbs sampling. Posterior samples can be obtained via Gibbs sampling. Sampling from the above hierarchy is expensive, as the full conditionals of λ_j and τ do not follow standard distributions and slice sampling has to be used. Makalic and Schmidt [14] propose an alternative Horseshoe hierarchy that exploits the following mixture representation of a half-Cauchy distributed random variable $X \sim \mathcal{C}^+(0, \Psi)$

$$(3.9) \quad X^2 | \psi \sim \mathcal{IG}\left(\frac{1}{2}, \frac{1}{\psi}\right),$$

$$(3.10) \quad \psi \sim \mathcal{IG}\left(\frac{1}{2}, \frac{1}{\Psi^2}\right),$$

which leads to conjugate conditional posterior distributions. The sampling scheme in [14] samples iteratively from the following set of full conditional posteriors

$$\begin{aligned} \beta | \cdot &\sim \mathcal{N}_p(\mathbf{A}^{-1} \mathbf{X}^T \mathbf{y}, \sigma^2 \mathbf{A}^{-1}) \\ \sigma^2 | \cdot &\sim \mathcal{IG}\left(\frac{n+p}{2}, \frac{(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)}{2} + \frac{\beta^T \Lambda_*^{-1} \beta}{2}\right) \\ \lambda_j^2 | \cdot &\sim \mathcal{IG}\left(1, \frac{1}{\nu_j} + \frac{\beta_j^2}{2\tau^2\sigma^2}\right) \\ \tau^2 | \cdot &\sim \mathcal{IG}\left(\frac{p+1}{2}, \frac{1}{\rho} + \frac{1}{2\sigma^2} \sum_{j=1}^p \frac{\beta_j^2}{\lambda_j^2}\right) \end{aligned}$$

$$v_j | \cdot \sim \mathcal{IG}\left(1, \frac{1}{A^2} + \frac{1}{\lambda_j^2}\right)$$

$$\rho | \cdot \sim \mathcal{IG}\left(1, 1 + \frac{1}{\tau^2}\right),$$

with $\mathbf{A} = (\mathbf{X}^T \mathbf{X} + \mathbf{\Lambda}_*^{-1})$, $\mathbf{\Lambda}_* = \tau^2 \mathbf{\Lambda}$, $\mathbf{\Lambda} = \text{diag}(\lambda_1^2, \dots, \lambda_p^2)$.

This sampling scheme allows computationally efficient sampling of the Horseshoe regularization. Sampling 1000 draws from the posterior distribution in the HorseRule model for the Boston housing data used in Section (4.4) takes about one minute. The complexity of the Horseshoe sampling depends mostly on the number of linear terms and decision rules, and increases only slowly with N . [13] suggest a computational shortcut where a given β_j is sampled in a given iteration only if the corresponding scale $(\lambda_j \cdot \tau)$ is higher than a threshold. The λ_j needs to be sampled in every iteration to give every covariate the chance of being chosen in the next iteration. We have implemented this approach and seen that it can give tremendous computational gains, but we have not used it when generating the results here since the effects it has on the invariant distribution of the MCMC scheme needs to be explored further.

3.5. Sampling the splitting points. The BART model can be seen as the sum of trees with a Gaussian prior on the terminal node values

$$\mu_j \sim \mathcal{N}\left(0, \frac{0.5}{\tau \sqrt{k}}\right),$$

where k denotes the number of trees. BART uses a fixed regularization parameter τ and samples the tree structure, while HorseRule uses a fixed rule structure and adapts to the data through sampling the shrinkage parameters λ_j and τ . Using a fixed tree structure offers dramatic computational advantages, as no Metropolis-Hastings updating steps are necessary, but the splits are likely to be suboptimal with respect to the whole ensemble.

As shown in Section (4), both HorseRule and BART achieve great predictive performance through different means, and a combination in which both shrinkage and tree structure are sampled in a fully Bayesian way could be very powerful, but computationally very demanding. An intermediate position is to keep the splitting variables fixed in HorseRule, but to sample the splitting points. We have observed that HorseRule often keeps very similar rules with slightly different splitting points in the ensemble, which is a discrete approximation to sampling the splitting points. Hence this could also improve interpretability since a large number of rules with nearby splitting points can be replaced by a single rule with an estimated splitting point. We leave this extension to future research, however.

4. EMPIRICAL RESULTS

4.1. Simulation Study. One advantage of RuleFit is the ability to complement the rules with linear terms or any other specified effect. This subsection analyses the ability of HorseRule and RuleFit to recover the true signal when the true relationship is linear and observed with noise. The data is generated with $X_i \sim \mathcal{N}(0, 1)$, $i = 1, \dots, 100$, $Y = 5X_1 + 3X_2 + X_3 + X_4 + X_5 + \epsilon$ and $\epsilon \sim \mathcal{N}(0, 1)$. The first 5 predictors thus have a positive dependency with y of

	RMSE			$\Delta\beta_{true}$			$\Delta\beta_{noise}$		
	$n = 100$	$n = 500$	$n = 1000$	$n = 100$	$n = 500$	$n = 1000$	$n = 100$	$n = 500$	$n = 1000$
OLS	3.23	1.10	1.06	1.25	0.19	0.14	2302	3.78	2.54
Horseshoe Regression	1.14	1.01	1.01	0.40	0.18	0.13	1.72	0.70	0.49
HorseRule $\alpha = 0, \beta = 0$	1.54	1.02	1.01	1.99	0.39	0.29	2.74	0.22	0.15
HorseRule $\alpha = 1, \beta = 2$	1.25	1.02	1.01	1.15	0.37	0.28	3.14	0.37	0.24
RuleFit $k = 2000$	1.84	1.23	1.15	3.58	1.42	1.05	1.18	0.91	0.99

TABLE 4.1. Simulation study. The true effect is linear.

varying magnitude while the remaining 95 covariates are noise. Table 4.1 report the results from 100 simulated datasets. RMSE measures the discrepancy between the fitted values and the true mean. RuleFit and HorseRule model use 500 rules in addition to the linear terms. The best model in RMSE is as expected the Horseshoe regression without any rules. The OLS estimates without any regularization struggles to avoid overfitting with all the unnecessary covariates and does clearly worse than the other methods. HorseRule without the rule structure prior outperforms RuleFit, but adding a rule structured prior gives an even better result. The differences between the models diminishes quickly with the sample size (since the data is rather clean), the exception being RuleFit which improves at a much lower rate than the other methods. Table 4.1 also breaks down the results into the ability to recover the true linear signal, measured by $\Delta\beta_{true} = |(\beta_1, \beta_2, \beta_3, \beta_4, \beta_5) - (5, 3, 1, 1, 1)|_1$, and the ability to remove the noise covariates, measured by $\Delta\beta_{noise} = |(\beta_6, \dots, \beta_{100}) - (0, \dots, 0)|_1$. We see that the HorseRule’s horseshoe prior is much better at recovering the true linear signal compared to RuleFit with its L1-regularization. OLS suffers from its inability to shrink away the noise.

Even though such clear linear effects are rare in actual application, the simulation results in Table 4.1 shows convincingly that HorseRule will prioritize and accurately estimate linear terms when they fit the data well. This is in contrast to RuleFit which shrinks the linear terms too harshly and compensates the lack of fit with many rules. HorseRule will only try to add non-linear effects through decision rules if they are really needed.

4.2. Influence of the rule generating process. In this section we will analyse the influence of different rule generating processes on model performance for the Diamond dataset with ($N = 308$ and $p = 4$) and the Boston housing data ($N = 506$ and $p = 13$). In each setting 1000 trees with an average tree depth of $L = 5$ are used, but the ensemble is generated by i) random forest, ii) gradient boosting or iii) a combination of 30% of the trees from random forest and 70% from gradient boosting.

Figure 4.1 shows the RMSE distribution over the folds used in 10-fold cross-validation. As expected the error-correcting rules found by gradient boosting outperforms randomly generated rules from Random Forest on both datasets. However, combining the two types of rules leads to a lower mean RMSE on both datasets, and much lower median RMSE in the Boston dataset. Moreover, the RMSE is less variable over the 10 folds when using a mixture of ensembles.

In our experiments it rarely hurts the performance to use both type of rules, but on some datasets it leads to a dramatically better prediction accuracy. The mixing proportion for

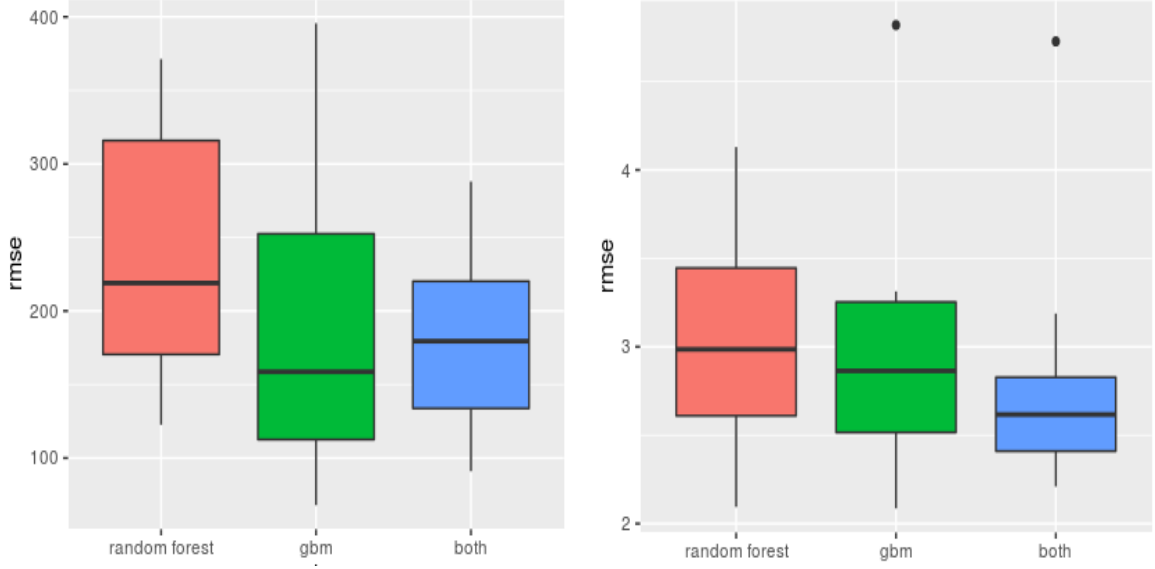


FIGURE 4.1. RMSE on the Diamonds (left) and the Boston (right) dataset for different rule generation strategies.

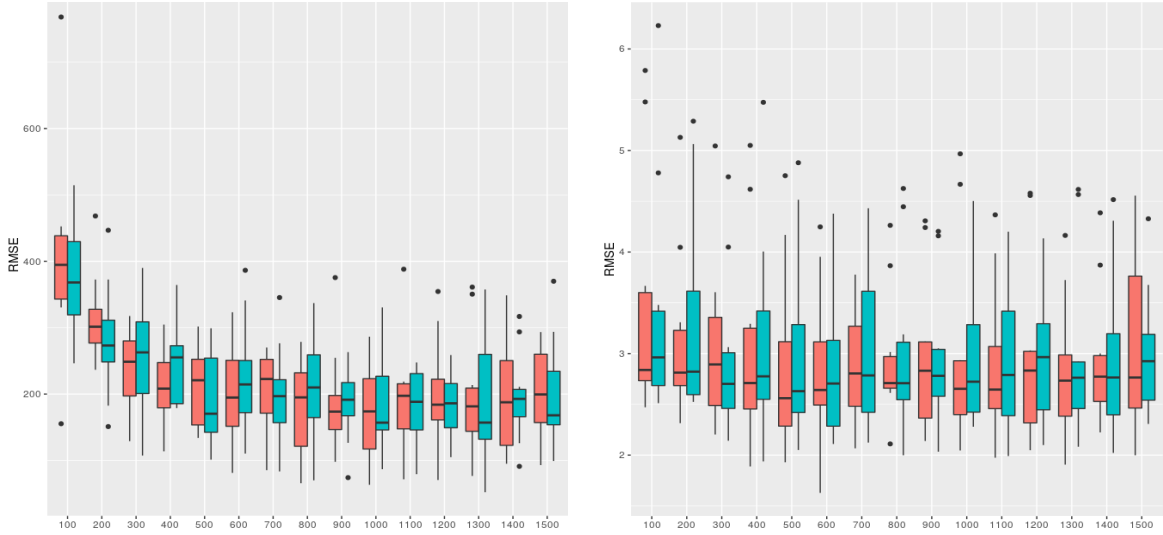


FIGURE 4.2. RMSE depending on the number of trees on the diamonds (left) and Boston (right) dataset for $(\mu, \eta) = (0, 0)$ (red) and $(\mu, \eta) = (1, 2)$ (blue).

the ensemble methods can also be seen as a tuning parameter to give a further boost in performance.

Another parameter that is potentially crucial is the number of trees used to generate the decision rules. In gradient boosting limiting the number of trees (iterations) is the most common way to control overfitting. Also in BART the number of trees has a major impact on the quality and performance of the resulting ensemble [4]. The same is expected for RuleFit, as it uses L1-regularization; with an increasing number of rules the overall shrinkage λ increases, leading to an over-shrinkage of good rules.

TABLE 4.2. Settings for the compared methods.

Method	Parametersetting
HR-default	$L = 5$ Ensemble: GBM+RF $(\mu, \eta) = (1, 2)$
HR-CV	$L = (2, 5, 8)$ Ensemble: GBM+RF $(\mu, \eta) = ((0, 0), (0.5, 0.5), (1, 2))$
RuleFit	$k = 500, 1000, \dots, 5000$ $L = (2, 5, 8)$
Random Forest	% Variables tried = $(0.25, 0.5, 0.75, 1, \sqrt{p})$
BART	$(\gamma, q) = ((3, 0.9), (3, 0.99), (10, 0.75))$ $\tau = 2, 3, 5$ Number of Trees: 50,200

TABLE 4.3. The 16 regression datasets.

Name	N	Q	C
Abalone	4177	7	1
Ais	202	11	1
Attend	838	6	3
Baskball	96	4	0
Boston	506	13	0
Budget	1729	10	0
Cps	534	7	3
Cpu	209	6	1
Diamond	308	1	3
Hatco	100	6	4
Heart	200	13	3
Fat	252	14	0
Mpg	392	6	1
Ozone	330	8	0
Servo	167	2	2
Strike	625	4	1

To investigate the sensitivity of HorseRule to the number of trees, we increase the number of trees successively from 100 to 1500 in the Boston and Diamonds datasets. This corresponds to $500, 550, \dots, 5 \cdot 1500 = 7500$ decision rules before removing duplicates. We also test if the rule structured prior interacts with the effect of the number of trees by running the model with $(\mu, \eta) = (0, 0)$ and $(\mu, \eta) = (1, 2)$. Figure 4.2 shows the performance of HorseRule as a function of the number of trees used to extract the rules. Both HorseRule models are relatively insensitive to the choice of k , unless the number of trees is very small. Importantly, no overfitting effect can be observed, even when using an extremely large number of 1500 trees on relatively small datasets ($N = 308$ and $N = 506$ observations, respectively). We will use 1000 trees as a standard choice in the following, but a small number of trees can be used if computational complexity is an issue, with little to no expected loss in accuracy.

4.3. Prediction performance comparison on 16 datasets. We compare the predictive performance of HorseRule with competing methods on 16 regression datasets. The datasets are the subset of the datasets used in [4] that were available to us online. Also datasets with clearly non-gaussian response variable even after transformations are excluded. We use 10-fold cross validation on each data set and report the relative RMSE (RRMSE) in each fold; RRMSE for a fold is the RMSE for a method divided by the RMSE of the best method on that fold. This allows us to compare performance over different datasets with differing scales and problem difficulty.

To ensure a fair comparison we use another (nested) 5-fold cross validation in each fold to find good values of the tuning parameters for each method. For BART and Random Forest the cross-validation settings from [4] are used. For RuleFit we cross-validate over the number of rules and the depth of the trees, as those are the potentially most impactful parameters. The shrinkage τ in RuleFit is determined by the model internally. For HorseRule we use cross-validation to identify suitable hyperparameters (μ, η) as well as the tree depth. We also run a HorseRule version with the proposed standard settings without cross-validation. Table 4.2 summarizes the settings of all methods and Table 4.3 the characteristics of the datasets.

Figure 4.3 shows a comparison of the different HorseRule models over $10 \cdot 16 = 160$ data cross-validation splits. While the $(\mu, \eta) = (1, 2)$ already performs better than the prior without rule structure $((\mu, \eta) = (0, 0))$, cross-validation of (μ, η) helps to improve performance further.

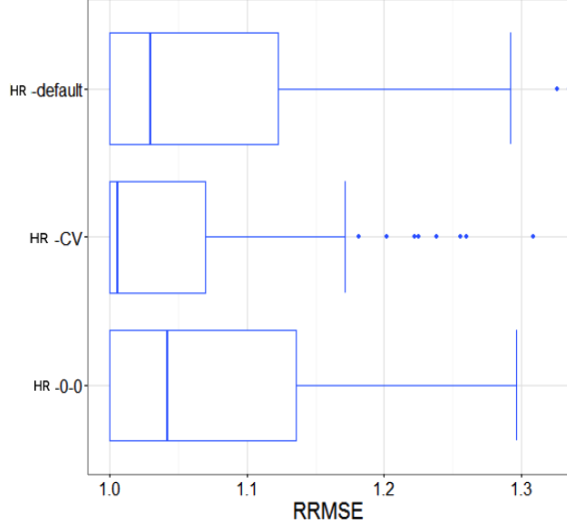


FIGURE 4.3. Relative RMSE comparison of the different HorseRule versions.

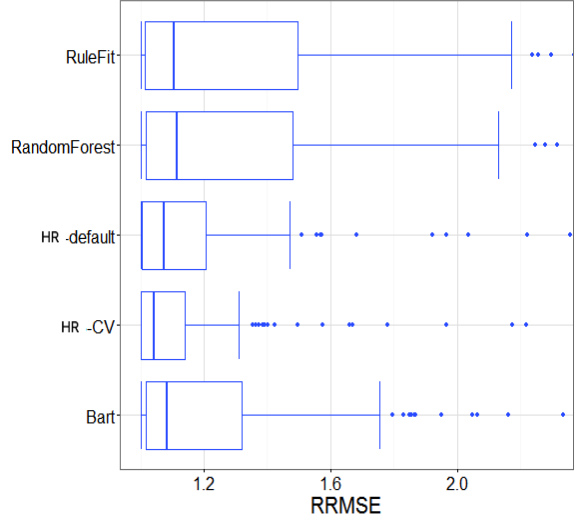


FIGURE 4.4. Relative RMSE comparison of HorseRule with competing methods.

Figure 4.4 shows that HorseRule has very good performance across all datasets and folds, and the median RRMSE is substantially smaller than its competitors. Table 4.4 shows that HorseRule-CV is the best model on 6/16 datasets and has the lowest rank. Moreover, the last row of Table 4.4 displays the worst RRMSE over all datasets for each method; it shows that whenever HorseRule-CV is not the best method, it is only marginally behind the winner. This is not true for BART, RandomForest and RuleFit, which all perform substantially worse than the best method on some datasets. The small differences between the default settings for HorseRule and HorseRule-CV shows that HorseRule performs well also without time-consuming tuning of its parameters. Note also that HorseRule without cross-validation performs better on average than the competing methods with cross-validated parameters.

It is interesting to note that HorseRule tends to perform well on datasets in which BART does not perform well and vice versa, despite their similar mindset. RuleFit performs the best on 2/16 datasets, and the median RRMSE is slightly lower than for RandomForest.

TABLE 4.4. Cross-validated prediction performance for the 16 regression data sets. Each entry shows the RMSE and in parentheses the rank on this data set. The best result is marked in bold.

	BART	RandomForest	RuleFit	HorseRule	HorseRule-CV
Abalone	2.150 (5)	2.119 (3)	2.139 (4)	2.115(2)	2.114 (1)
AIS	1.144 (3)	1.247 (5)	1.207 (4)	0.713 (2)	0.699 (1)
Attend	394141 (3)	411900 (5)	345177 (1)	398485 (4)	365010 (2)
Baskball	0.087 (2)	0.086 (1)	0.088 (3.5)	0.088 (3.5)	0.092 (5)
Boston	2.867 (1)	3.153 (5)	3.037 (4)	2.940 (3)	2.926 (2)
Budget	0.039 (2)	0.038 (1)	0.061 (5)	0.041 (3)	0.042 (4)
CPS	4.356 (2)	4.399 (5)	4.386 (4)	4.348 (1)	4.370 (3)
Cpu	41.52 (3)	54.08 (4)	54.50 (5)	36.03 (1)	37.47 (2)
Diamond	215.0 (3)	465.9 (4)	233.7 (5)	184.5 (2)	171.27 (1)
Hacto	0.453 (5)	0.311 (4)	0.297 (3)	0.261 (2)	0.260 (1)
Heart	8.917 (1)	9.048 (2)	9.349 (5)	9.241 (4)	9.070 (3)
Fat	1.306 (5)	1.114 (1)	1.173 (2)	1.264 (4)	1.245 (3)
MPG	2.678 (2)	2.692 (4)	2.672 (1)	2.714 (5)	2.689 (3)
Ozone	4.074 (2)	4.061 (1)	4.189 (5)	4.120 (3)	4.165 (4)
Servo	0.588 (5)	0.486 (3)	0.502 (4)	0.409 (2)	0.403 (1)
Strikes	458.4 (5)	453.7 (4)	447.7 (2)	449.2 (3)	447.2 (1)
Average Rank	3.062	3.188	3.719	2.656	2.375
Worst Relative RRMSE	1.742	2.720	1.726	1.154	1.117

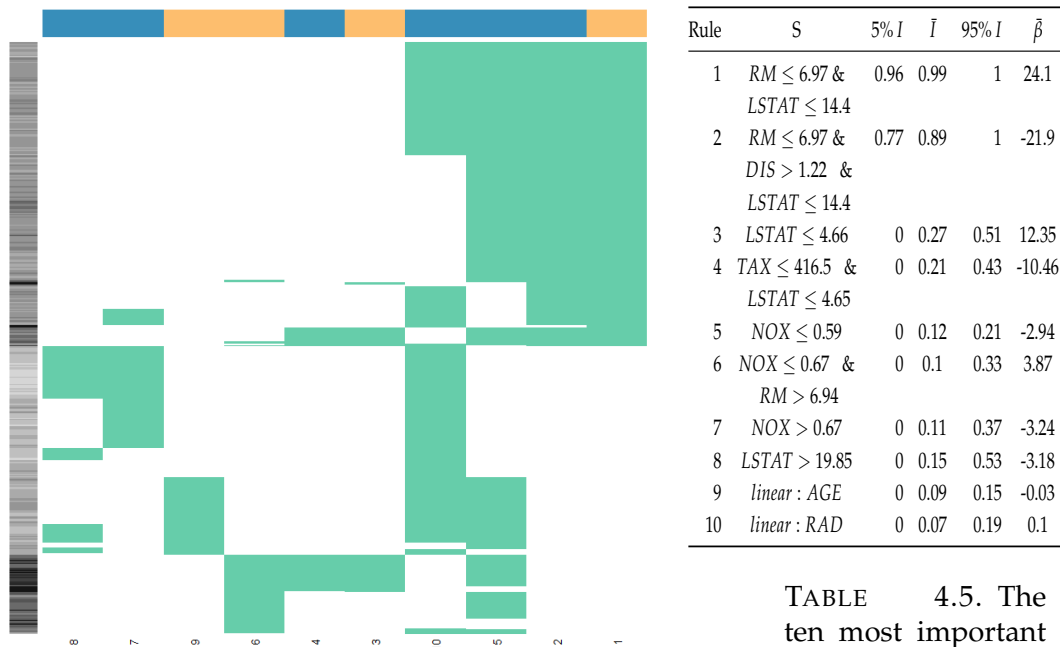


FIGURE 4.5. RuleHeat for the Boston Housing data. See the text for details.

TABLE 4.5. The ten most important rules in the Boston housing data

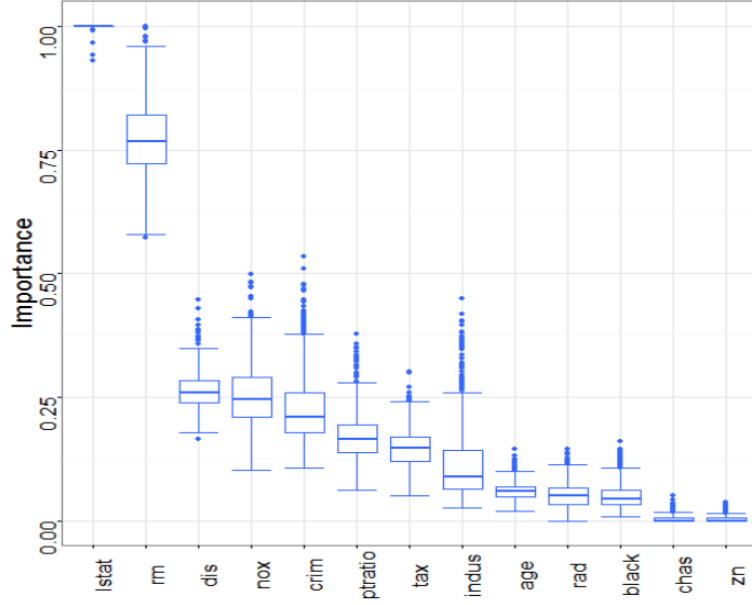


FIGURE 4.6. Posterior distribution of the Input Variable Importance for the 13 Covariates.

4.4. Boston Housing. In this section we apply HorseRule to the well known Boston Housing dataset to showcase its usefulness in getting insights from the data. The Boston housing data consists of $N = 506$ observations which are city areas in Boston and $p = 13$ covariates variables are recorded. These variables include ecological measures of nitrogen oxides (NOX), particulate concentrations (PART) and proximity to the Charles River (CHAS), the socio-economic variables proportion of black population (B), property tax rate (TAX), proportion of lower status population (LSTAT), crime rate (CRIM), pupil teacher ratio (PTRATIO), proportion of old buildings (AGE), the average number of rooms (RM), area proportion zoned with large lots (ZN), the weighted distance to the employment centers (DIS) and an index of accessibility to key infrastructure (RAD). The dependent variable is the median housing value in the area. The HorseRule with default parameter settings is used to fit the model. Table 4.5 shows the 10 most important effects. Following [9], the importance of a linear term is defined as

$$I(x_j) = |\beta_j| sd(x_j)$$

where $sd(\cdot)$ is the standard deviation, and similarly for a predictor from a decision rule

$$I(r_l) = |\alpha_l| sd(r_l).$$

We use the notation I_j when it is not important to distinguish between a linear term and a decision rule. For better interpretability we normalize the importance to be in $[0, 1]$, so that the most important predictor has an importance of 1. Table 6 reports the posterior distribution of the normalized importance (obtained from the MCMC draws) of the 10 most important rules or linear terms. The most important single variable is LSTAT, which appears in many of the rules, and as a single variable in the third most important rule. Note also that LSTAT does not appear as a linear predictor among the most important predictors.

To interpret the more complex decision rules in Table 4.5 it is important to understand that decision rules in an ensemble have to be interpreted with respect to other decision rules, and in relation to the data points covered by a rule. A useful way to explore the effects of the most important rules is what we call a *RuleHeat* plot, see Figure 4.5 for an example for the Boston housing data. The horizontal axis lists the most important decision rules and the vertical axis the N observations. A square is green if $r_l(\mathbf{x}) = 1$. The grayscale on the bar to the left indicates the outcome (darker for higher price) and the colorbar in the top of the figure indicates the sign of the covariate’s coefficient in the model (sand for positive). RuleHeat makes it relatively easy to find groups of similar observations, based on the rules found in HorseRule, and to assess the role a rule plays in the ensemble. For example, Figure 4.5 shows that the two most important rules differ only in a few observations. The two rules have very large coefficients with opposite signs. Rule 1 in isolation implies that prices are substantially higher when the proportion of lower status population is low ($LSTAT \leq 14.4$) for all but the very largest houses ($RM \leq 6.97$). However, adding Rule 2 essentially wipes out the effect of Rule 1 ($24.1 - 21.9 = 2.2$) except for the six houses very close to the employment centers ($DIS < 1.22$) where the effect on the price remains high.

Similarly to the Variable importance in Random Forest and RuleFit, we can calculate a variable input importance for the HorseRule model. The importance of the j th predictor given the data is defined as [9]

$$J(x_j) = I(x_j) + \sum_{l: j \in Q_l} I(r_l) / |Q_l|$$

where the sum runs over all rules where x_j is one of the predictors use to define the rule. Note how the importance of the rules are discounted by the number of variables involved in the rule, $|Q_l|$. Figure 4.6 shows

the posterior distribution of $J(x_j)$ for the 13 covariates. LSTAT is the most important covariate with median posterior probability of 1 and very narrow posterior spread, followed by RM which has a median posterior probability of around 0.75. The importance of some variables, like NOX and INDUS, has substantial posterior uncertainty whereas for other covariates, such as AGE, the model is quite certain that the importance is low.

The overlapping rules, as well as similar rules left in the ensemble in order to capture model uncertainty about the splitting points make interpretation somewhat difficult. One way to simplify the output from HorseRule is to use the *decoupling shrinkage and summary* (DSS) approach by Carvalho and Hahn [12]. The idea is to reconstruct the full posterior estimator $\hat{\beta}$ with a 1-norm penalized representation, that sets many of the coefficients to exactly zero and also merges together highly correlated coefficients. We do not report systematical tests here, but in our experiments using DSS with a suitable shrinkage parameter did not hurt the predictive performance, while allowing to set a vast amount of coefficients to zero. Using HorseRule followed by DSS on the Boston housing data leaves 106 non-zero coefficients in the ensemble. The 10 most important rules can be seen in table 4.6. We can see that the new coefficients are now less overlapping. The relatively small number of rules simplify interpretation. Posterior summary for regression with shrinkage priors is an active field of

TABLE 4.6. The ten most important rules in Boston data after DSS.

Rule	S	\bar{Imp}	$\bar{\beta}$
1	$RM \leq 7.13$	1	-3.47
2	$RM \leq 6.98$ & $PTRATIO \leq 18.7$ & $LSTAT > 5.95$	0.97	-2.36
3	$LSTAT > 18.75$	0.81	1.80
4	$linear : RAD$	0.80	0.10
5	$RM \leq 7.437$ & $LSTAT \leq 7.81$	0.79	-2.03
6	$NOX \leq 0.62$ & $RM \leq 7.31$	0.70	-1.64
7	$RM \leq 7.1$ & $RAD \leq 4.5$ & $LSTAT \leq 7.81$	0.68	-2.47
8	$NOX > 0.59$	0.63	-1.47
9	$linear : LSTAT$	0.58	-0.09
10	$linear : AGE$	0.58	-0.02

research (see e.g. [15]) and future developments might help to simplify the rule ensemble further.

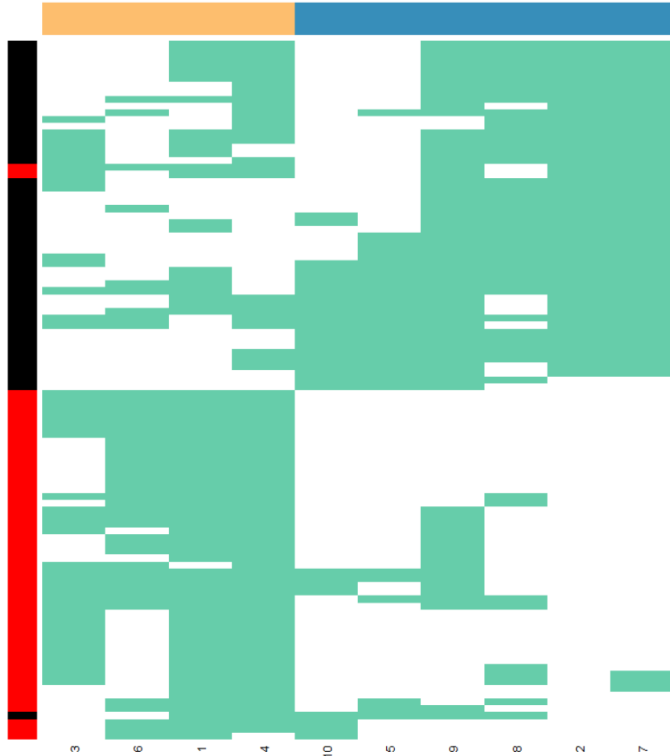


FIGURE 4.7. RuleHeat for the prostate cancer data. Details can be found in the text. Cancer Samples are colored in red, normal in black.

	BART	RandomForest	RuleFit	HorseRule
CV-Accuracy	0.900	0.911	0.831	0.922
CV-AUC	0.923	0.949	0.953	0.976
Test-Accuracy	0.824	0.971	0.941	0.971
Test-AUC	1	0.991	0.995	1

TABLE 4.7. Cross-validated accuracy on Crossvalidation and the test data.

Rule	S	5%	\bar{Imp}	95%	$\bar{\beta}$
1	$556_s_at \leq 1.55$	0	0.33	1	3.10
2	$34647_at \leq -1.18$ & $37639_at \leq 1$	0	0.15	1	-1.78
3	$37478 > -0.32$	0	0.18	0.91	1.42
4	$38087_s_at \leq 0.83$	0	0.23	1	1.81
5	$34678_at > 0.38$	0	0.19	0.88	-1.58
6	$1243_at \leq 0.35$	0	0.15	0.66	1.19
7	$37639_at \leq 1$	0	0.13	0.80	-1.10
8	$33121_g_at \leq 0.672$ & $960_g_at > 0.378$	0	0.10	0.82	-1.09
9	$41706_at \leq 1.33$	0	0.15	0.79	-1.13
10	$39061_at > 0.31$	0	0.1	0.52	-1.03

TABLE 4.8. The ten most important rules found in the prostate cancer data

4.5. Cancer Data. In this section we analyze how HorseRule can be used to find interesting pattern in a classification problems, for example using logistic regression on gene expression data to find genes that can signal the presence or absence of cancer. This is an area where interpretability clearly matters. Extending the regression model to logistic regression can be easily done with the introduction of a latent variable. We chose to use the Pólya–Gamma latent variable scheme by [16]. Methodological difficulties arise from the usually small number of available samples, as well as high number of candidate genes, leading to an extreme $p \gg n$ situation. We showcase the ability of HorseRule to make inference in this difficult domain on the Prostate Cancer dataset, which consists of 52 cancerous and 50 normal samples ($n = 102$). In the original data $p = 12600$ genetic expressions are available, which can be reduced to 5966 genes after applying the pre-processing described in Singh et al. [18]. Since spurious relationships can easily occur when using higher order interactions in the $p \gg n$ situation, we use the hyperparameters $\mu = 2$ and $\eta = 4$ to express our prior belief, that higher order interactions are very unlikely to be reflect any true mechanism.

Table 4.7 shows that HorseRule has higher accuracy and significantly higher AUC than the competing methods. We also test the methods on a unseen test data set containing 34 samples not used in the previous step. All methods have lower error here, implying that the test data consists of more predictable cases. The difference is smaller, but HorseRule performs slightly better here as well.

The 10 most important rules for HorseRule are found in Table 4.8. It contains 8 rules with one condition and only 2 with two conditions, implying that there is not enough evidence in the data for complicated rules to overrule our prior specification. All of the most important rules still contain 0 in their 5% posterior importance distribution, implying that they are eliminated by the model in at least 5% of the samples; the small sample size leads to non-conclusive results.

Figure 4.7 shows the RuleHeat plot for the 10 most important rules. The outcome is binary, and the vertical bar to the left is red for cancer and black for normal. The most important Rule 1 covers all except one cancer patient. This gene would probably not be found to be significant using traditional tests in logistic regression, as it is individually not very discriminative. Its importance arises from the combination with the other rules, especially Rule 2, Rule 7 and Rule 8, that are able to correct the false positive predictions using Rule 1 alone.

Figure 4.8 and Figure 4.9 show the subspaces from the two most important interaction rules. Again normal samples are colored black and cancerous red. The first interaction looks somewhat unnatural. A strong effect can be assumed for *37639_at* where higher values indicate cancer. This rule is also individually represented as Rule 7. The second split on *34647_at* < -1.18 corrects 3 missclassified observations by the first split alone. This rule probably only works well in the ensemble but may not reflect a true mechanism.

The second interaction effect is more interesting. It seems that normal samples have lower values in *33121_g_at* and higher values in *960_g_at*. This rule might reflect a true interaction mechanism and could be worth analysing further.

Figure 4.10 shows the input variable importance of the 50 most important genes. In this domain the advantage of having estimates of uncertainty can be very beneficial, as biological

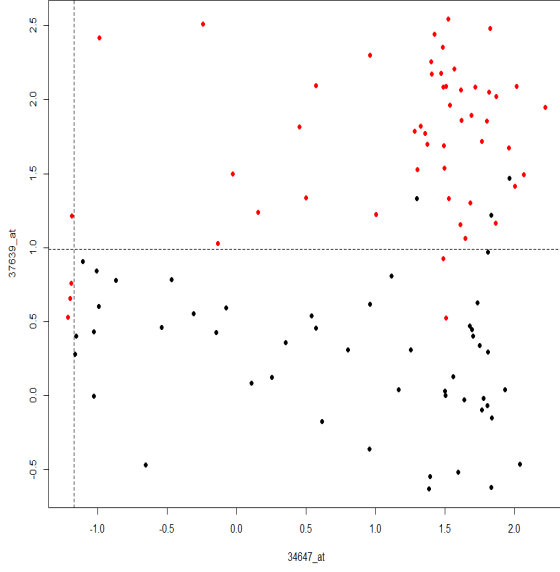


FIGURE 4.8. Scatterplot for Genes 37639_at and 34647_at. Normal samples in black and cancerous samples in red. Rule 2 is defined by the bottom right quadrant.

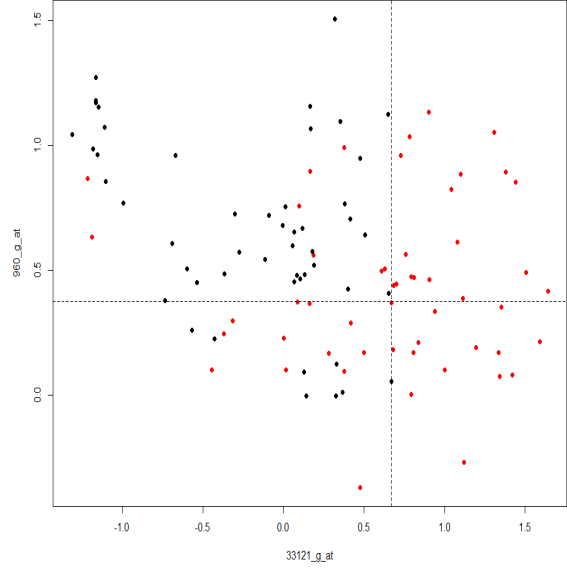


FIGURE 4.9. Scatterplot for Genes 33121_g_at and 960_g_at. Normal samples in black and cancerous samples in red. Rule 8 is defined by the top left quadrant.

follow up studies are costly and the probability of spurious relationships is high. In this data the genes 37639_at and 556_s_at contain an importance of 1 in their 75 % confidence bands, and are therefore with high certainty the most important genes for the prediction of prostate cancer. 37639_at was found in previous studies to be associated with prostate cancer, however 556_s_at seems to be understudied in this context but was already found to be important for breast cancer. This gene is represented in Rule 1. As discussed above the gene is individually not very discriminative, but only in conjunction with other rules.

CONCLUSIONS

We propose HorseRule, a new model for flexible non-linear regression and classification. The model is based on RuleFit and uses decision rules from a tree ensemble as predictors in a regularized linear fit. We replace the L1-regularization in RuleFit with a horseshoe prior with a hierarchical structure especially tailored for a situation with decision rules as predictors. Our prior shrinks complex (many splits) and specific (small number of observations satisfy the rule) rules more heavily a priori, and is shown to be efficient in removing noise without tampering with the signal. The efficient shrinkage properties of the new prior also makes it possible to complement the rules from boosting used in RuleFit with an additional set of rules from random forest. The rules from random forest are not as tightly coupled as the ones from boosting, and are shown to improve prediction performance compared to using only rules from boosting.

HorseRule is demonstrated to outperform state-of-the-art competitors like RuleFit, BART and random forest in an extensive evaluation of predictive performance on 16 widely used

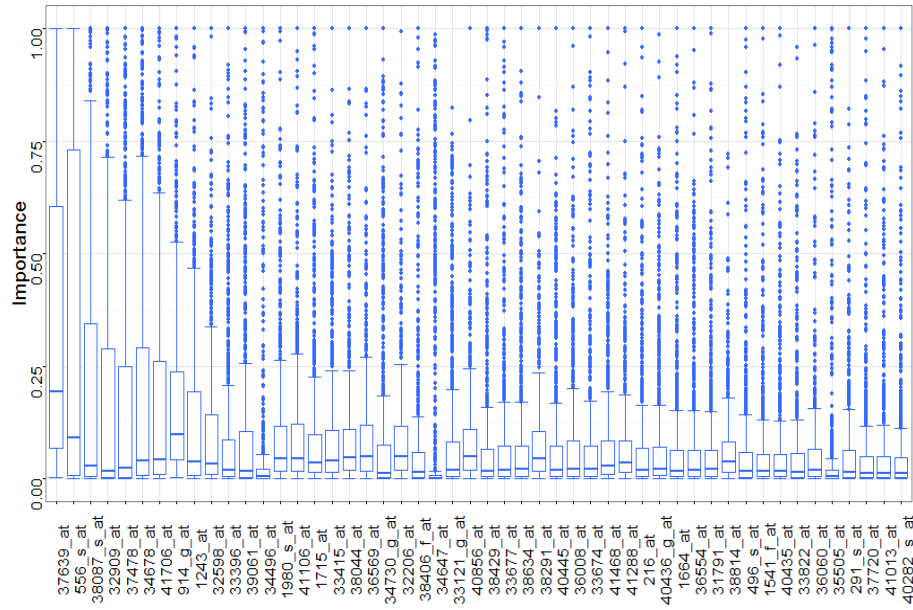


FIGURE 4.10. Posterior distribution of the Input Variable Importance of the 50 most influential Covariates.

datasets. Importantly, HorseRule performs consistently well on all datasets, whereas the other methods perform quite poorly on some of datasets. We also demonstrate the interpretation of HorseRule in both a regression and a classification problem. HorseRule's use of decision rules as predictors and its ability to keep only the important predictors makes it easy to interpret its results, and to explore the importance of individual rules and predictor variables.

APPENDIX A - THE HORSERULE R PACKAGE

The following code illustrates the basic features of ourHorseRule package in R with standard settings.

```
library(horserule)
data(Boston)

N = nrow(Boston)
train = sample(1:N, 500)
Xtrain = Boston[train,-14]
ytrain = Boston[train, 14]
Xtest = Boston[-train, -14]
ytest = Boston[-train, 14]

# Main function call (variable scaling performed internally)
hrres = HorseRuleFit(X=Xtrain, y=ytrain,
  # MCMC settings
  thin=1, niter=1000, burnin=100,
  # Parameters for the rule generation process
  L=5, S=6, ensemble = "both", mix=0.3, ntree=1000,
  # Model parameters. Data is scaled so no intercept needed.
  intercept=F, linterms=lin, ytransform = "log",
```

```

# Hyperparameters for the rule structured prior
alpha=1, beta=2, linp = 1, restricted = 0)

# Check model performance by predicting holdout cases
pred = predict(hrres, Xtest)
sqrt(mean((pred-ytest)^2))

# Find most important rules
importance_hs(hrres)

# Compute variable importance
variable_importance(hrres)

# Monitor the complexity of the rules
complexity_plot(hrres)

```

REFERENCES

- [1] Leo Breiman. Random forests. *Machine Learning*, 45:5–32, 2001.
- [2] Carlos M Carvalho, Nicholas G Polson, and James G Scott. Handling sparsity via the horseshoe. In *AISTATS*, volume 5, pages 73–80, 2009.
- [3] Carlos M. Carvalho, Nicholas G. Polson, and James G. Scott. The horseshoe estimator for sparse signals. *Biometrika*, 97(2):465–480, 2010.
- [4] Hugh A. Chipman, Edward I. George, and Robert E. McCulloch. BART: Bayesian additive regression trees. *Annals of Applied Statistics 2010, Vol. 4, No. 1*, 266–298, October 2010.
- [5] William Weston Cohen. Fast Effective Rule Induction. In *Proceedings of the 12th International Conference on Machine Learning (ICML’95)*, pages 115–123, 1995.
- [6] Krzysztof Dembczyński, Wojciech Kotłowski, and Roman Słowiński. Ender: a statistical framework for boosting decision rules. *Data Mining and Knowledge Discovery*, 21(1):52–90, 2010.
- [7] Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29:1189–1232, 2000.
- [8] Jerome H Friedman and Bogdan E Popescu. Importance sampled learning ensembles. *Journal of Machine Learning Research*, 94305, 2003.
- [9] Jerome H. Friedman and Bogdan E. Popescu. Predictive learning via rule ensembles. *Annals of Applied Statistics 2008, Vol. 2, No. 3*, 916–954, November 2008.
- [10] Johannes Fürnkranz. Separate-and-conquer rule learning. *Artificial Intelligence Review*, 13(1):3–54, 1999.
- [11] Edward I. George and Robert E. McCulloch. Variable Selection via Gibbs Sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- [12] P Richard Hahn and Carlos M Carvalho. Decoupling shrinkage and selection in bayesian linear models: a posterior summary perspective. *Journal of the American Statistical Association*, 110(509):435–448, 2015.

- [13] Longhai Li and Weixin Yao. Fully bayesian logistic regression with hyper-lasso priors for high-dimensional feature selection. *arXiv preprint arXiv:1405.3319*, 2014.
- [14] Enes Makalic and Daniel F. Schmidt. A Simple Sampler for the Horseshoe Estimator. *IEEE Signal Process. Lett.*, 23(1):179–182, 2016.
- [15] Juho Piironen and Aki Vehtari. Comparison of bayesian predictive methods for model selection. *Statistics and Computing*, pages 1–25, 2016.
- [16] Nicholas G Polson, James G Scott, and Jesse Windle. Bayesian inference for logistic models using pólya–gamma latent variables. *Journal of the American statistical Association*, 108(504):1339–1349, 2013.
- [17] Robert E. Schapire. A Brief Introduction to Boosting. In *IJCAI*, pages 1401–1406, 1999.
- [18] Dinesh Singh, Phillip G Febbo, Kenneth Ross, Donald G Jackson, Judith Manola, Christine Ladd, Pablo Tamayo, Andrew A Renshaw, Anthony V D’Amico, Jerome P Richie, et al. Gene expression correlates of clinical prostate cancer behavior. *Cancer cell*, 1(2): 203–209, 2002.
- [19] Michael Smith and Robert Kohn. Nonparametric regression using Bayesian variable selection. *Journal of Econometrics*, 75(2):317–343, December 1996.
- [20] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996. ISSN 0035-9246.
- [21] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.